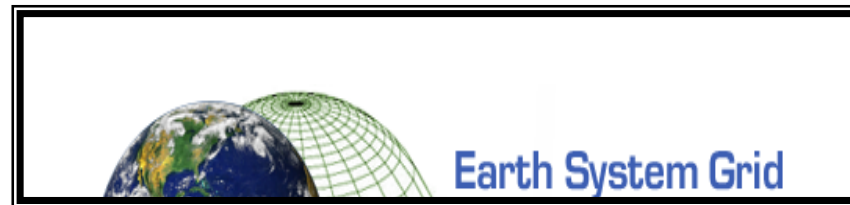# Earth System Grid: Model Data Distribution & Server-Side Analysis to Enable Intercomparison Projects
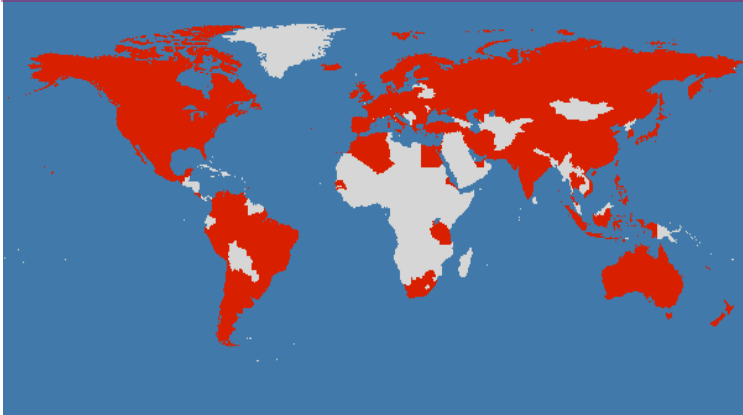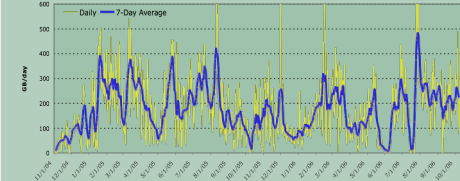


## PCMDI Software Team

# Challenges facing ESG-CET

- Building on the very successful CMIP3 IPCC AR4 ESG data portal.

- How best to collect and distribute data on a much larger scale?
  - At each stage tools could be developed to improve efficiency
  - Substantially more ambitious community modeling projects (>~300 TBs) will require a distributed database

- Metadata describing extended modeling simulations (e.g., atmospheric aerosols and chemistry, carbon cycle, dynamic vegetation, etc.)

- How to make information understandable to end-users so that they can interpret the data correctly

- More users from WGI. (Possibly WGII and WGIII?)

- Client and Server-side analysis and visualization tools in a distributed environment (i.e., subsetting, concatenating, regridding, filtering, …)

- Testbed needed by late 2008 – early 2009
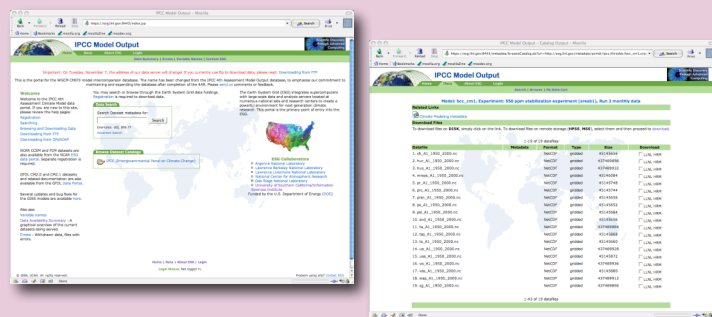
# ESG facts and figures

## ESG Objective

**To support the infrastructural needs of the national and international climate community, ESG is providing crucial technology to securely access, monitor, catalog, transport, and distribute data in today's Grid computing environment.**

## Worldwide ESG user base



## CMIP3 IPCC AR4 ESG Portal

**28 TB of data at the PCMDI site location**

- 68,400 files
- Generated by a modeling campaign coordinated by the Intergovernmental Panel on Climate Change
- Model data from 11 countries

**818 registered users**

**Downloads to date**

- 123 TB
- 543,500 files
- 300 GB/day (average)



IPCC Downloads (10/12/06)

Nov 2004 – Oct 2006

**200 scientific papers published to date based on analysis of CMIP3 IPCC AR4 data**

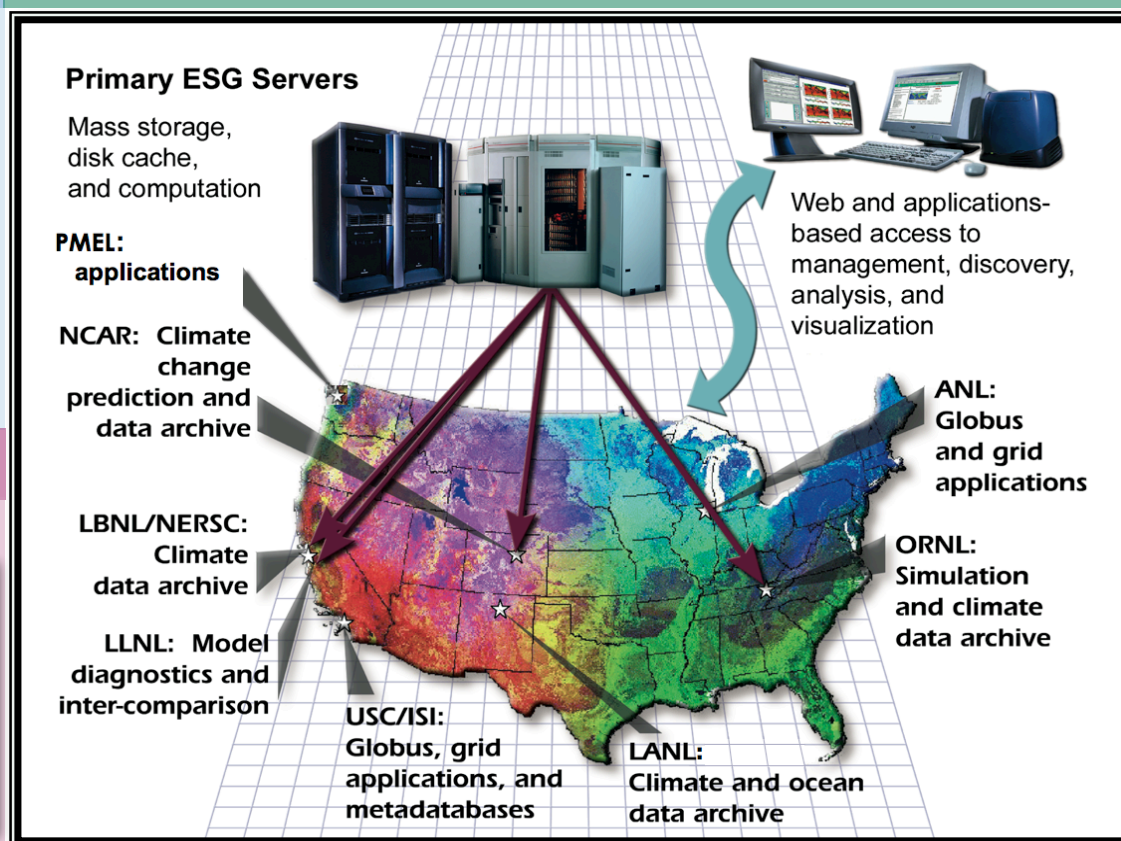# Providing climate scientists with virtual proximity to large simulation results needed for their research

Earth System Grid

## ESG Goal

- **Very large distributed data archives**
  - ➤ **Easy federation of sites**
  - ➤ **Across the US and around the world**
- **"Virtual Datasets" created through subsetting and aggregation**
- **Metadata-based search and discovery**
- **Web-based and analysis tool access**
- **Increased flexibility and robustness**
- **Server-side analysis**

### http://www-pcmdi.llnl.gov



## Current ESG Sites



**Primary ESG Servers**

Mass storage, disk cache, and computation

**PMEL:** applications

**NCAR:** Climate change prediction and data archive

**LBNL/NERSC:** Climate data archive

**LLNL:** Model diagnostics and inter-comparison

**USC/ISI:** Globus, grid applications, and metadatabases

Web and applications-based access to management, discovery, analysis, and visualization

**ANL:** Globus and grid applications

**ORNL:** Simulation and climate data archive

**LANL:** Climate and ocean data archive

# Evolving ESG for the future

Earth System Grid

## ESG Data System Evolution

### 2006

**Central database**

- **Centralized curated data archive**
- **Time aggregation**
- **Distribution by file transport**
- **No ESG responsibility for analysis**
- **Shopping-cart-oriented web portal**
- **ESG connection to desktop analysis tools (i.e., CDAT and CDAT-LAS)**

### Early 2009

**Testbed data sharing**

- **Federated metadata**
- **Federated portals**
- **Unified user interface**
- **Quick look server-side analysis with CDAT**
- **Location independence**
- **Distributed aggregation**
- **Manual data sharing**
- **Manual publishing**

### 2011

**Full data sharing** (add to testbed…)

- **Synchronized federation**
  - ➢ **metadata, data**
- **Full suite of server-side analysis with CDAT**
- **Model/observation integration**
- **ESG embedded into desktop productivity tools with CDAT**
- **GIS integration**
- **Model intercomparison metrics**
- **User support, life cycle maintenance**

**CCSM AR4**

## ESG Data Archive

**CCSM, AR5, satellite, In situ biogeochemistry, ecosystems**

**Terabytes**          **Petabytes**

# The growing importance of climate simulation data standards

- **Global Organization for Earth System Science Portal (GO-ESSP)**
  - ➤ **International collaboration to develop new generation of software infrastructure**
  - ➤ **Access to observed and simulated data from climate and weather communities**
  - ➤ **Working closely together using agreed upon standards**
  - ➤ **Last Annual meeting held at PCMDI**

- **NetCDF Climate and Forecast (CF) Metadata Convention standards**
  - ➤ **Specify syntax and vocabulary for climate and forecast metadata**
  - ➤ **Promotes the processing and sharing of data**
  - ➤ **The use of CF was essential for the success of the IPCC data dissemination**

# Supporting CF and CMOR

## Future issues for CF

- Develop further fundamental tools (such as Climate Model Output Rewriter - CMOR)
- Develop staggered and unstructured grids
- Deliver netCDF data into Geographical Information Systems (GIS)
- Upgrade to netCDF-4
- Include in situ observations

## CF/CMOR Development

- New CF website developed by PCMDI

- repository
  - News
  - Documents
    - ✓ CF Conventions
    - ✓ CF Standard Name table

- Conformance
  - Requirements & Recommendations
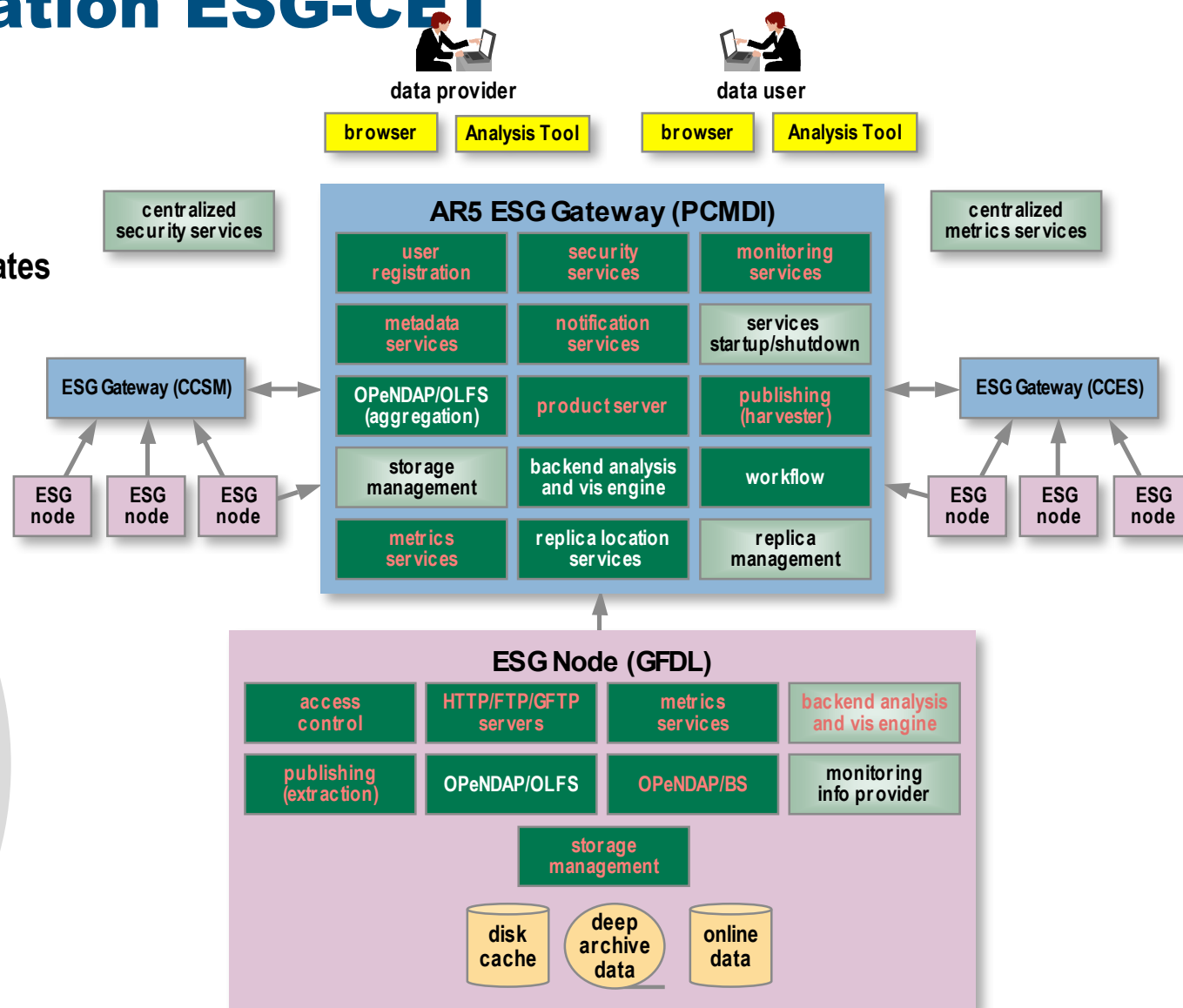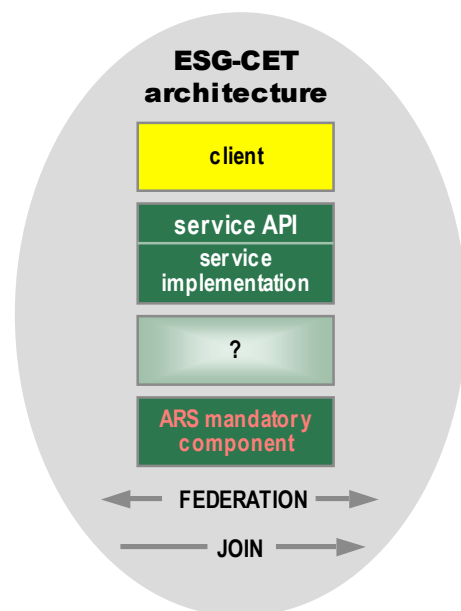  - CF Compliance Checker

- Mailing List
  - Archives

## New CF website

# Architecture of the next-generation ESG-CET

Earth System Grid

- **Huge data archives**
- **Broader geographical distribution of archives**
  - **across the United States**
  - **around the world**
- **Easy federation of sites**
- **Increased flexibility and robustness**

data provider

browser | Analysis Tool

data user

browser | Analysis Tool

centralized security services

centralized metrics services

**AR5 ESG Gateway (PCMDI)**

| user registration | security services | monitoring services |
| metadata services | notification services | services startup/shutdown |
| OPeNDAP/OLFS (aggregation) | product server | publishing (harvester) |
| storage management | backend analysis and vis engine | workflow |
| metrics services | replica location services | replica management |

ESG Gateway (CCSM)

ESG node | ESG node | ESG node

ESG Gateway (CCES)

ESG node | ESG node | ESG node

**ESG-CET architecture**

client

service API

service implementation

?

ARS mandatory component

← FEDERATION →

JOIN →

**ESG Node (GFDL)**

| access control | HTTP/FTP/GFTP servers | metrics services | backend analysis and vis engine |
| publishing (extraction) | OPeNDAP/OLFS | OPeNDAP/BS | monitoring info provider |

storage management

disk cache | deep archive data | online data

# Climate Data Analysis Tools: Software for Distributed Model Diagnosis & Intercomparison Research



## PCMDI Software Team
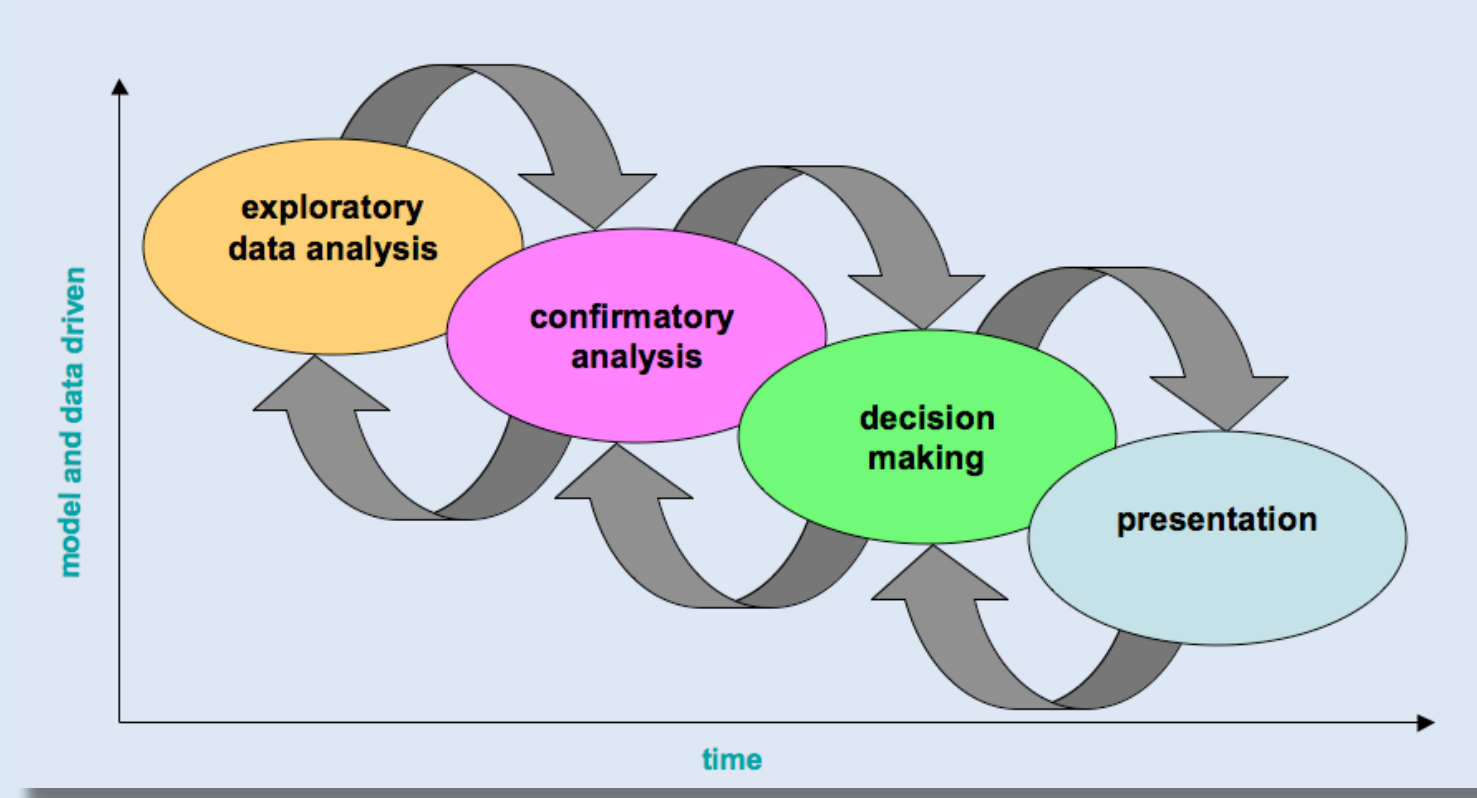
# Challenges facing CDAT

- **Integrating CDAT into a distributed environment**

- **Providing climate diagnostics**

- **Delivering climate component software to the community**

- Working with other forms of climate Metadata describing extended modeling simulations (e.g., atmospheric aerosols and chemistry, carbon cycle, dynamic vegetation, etc.)

- Testbed needed by late 2008 – early 2009

# CDAT objectives

## CDAT Objectives

**Seamless mechanisms for climate information exploration and analysis.**

# Enabling data management, data analysis, and visualization for intercomparison research
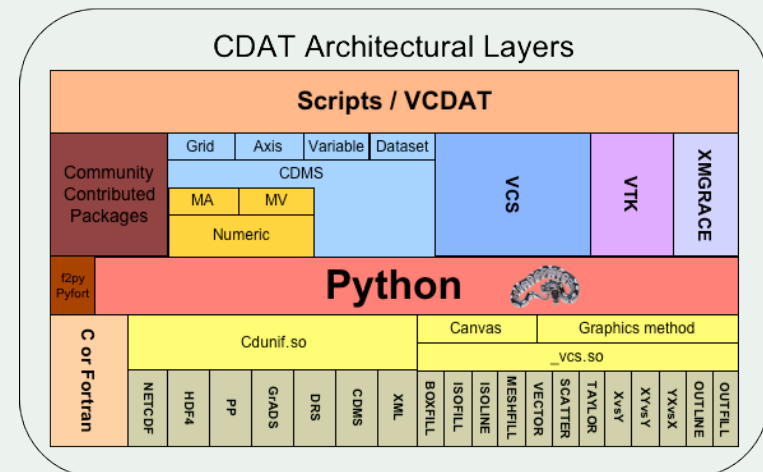
| CDAT Goal | What is CDAT? |
|---|---|
| Address the challenges of enabling data management, discovery, access, and advanced data analysis for climate model diagnosis and intercomparison research. | • CDAT *IS* Python!<br><br>• Designed for climate science data<br><br>• Scriptable<br><br>• Open-source and free |

## Typical usage examples of CDAT

- Calculate a long-term average

- Define wind-speed from u- and v-components

- Subset a dataset, selecting a spatiotemporal region

- Aggregate 1000s of files into a small XML file

- Generate a Hovmoller plot



CDAT Architectural Layers

# Evolving CDAT into an integrated client technology workplace

## CDAT Integrated Analysis Evolution

### 2006
**Community software**

- **Python based**
- **Start to finish environment**
- **Diverse analysis tools**
- **Languages: C/C++, Java, FORTRAN, Python**
- **Platforms: Unix, Mac, Windows**
- **VCDAT: discover, learn, and browse with a few clicks**
- **Connection to ESG**

### Early 2009
**Testbed distributed analysis**

- **Equal-access to shared resources (Web/Grid services)**
- **Quick look server-side analysis tool for ESG**
- **Diagnostics specific to AR5**
- **GFDL Ncvtk 3D visualization**
- **Web-CDAT: discover, learn, and browse via web browser**
- **Serving Google Maps and Google Earth Data with CDAT**

### 2011
**Full analysis sharing**

- **Full suite server-side analysis tool for ESG**
- **ESG embedded into desktop productivity tools (i.e., CDAT)**
- **GIS integration with CDAT**
- **SciDAC VACET analysis and visualization collaboration**
- **Global Organization for Earth System Science Portal (GO-ESSP)**
- **Remote generic apps for ESG**

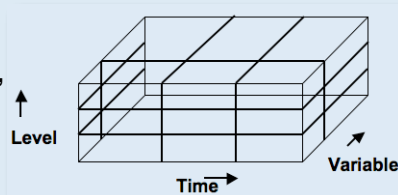**CDMS Numeric / MV Genutil / Cdutil VCS**

## CDAT Core Modules

**CDMS, Numeric, Genutil, Cdutil, Ncvtk, VACET, Diagnostics, ESG**

**Standalone** — **Distributed**

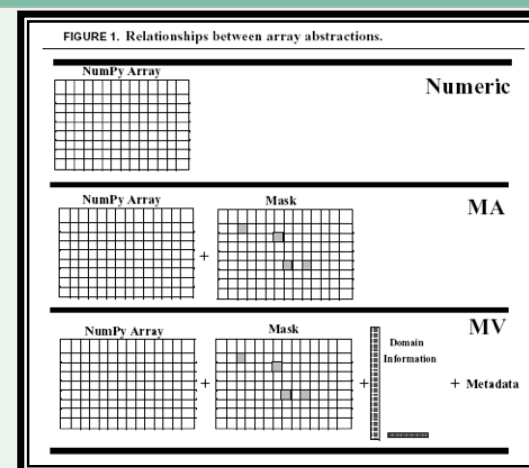# CDAT examples

| CDSCAN | MV |
|---|---|

**CDSCAN**

- Data aggregation: collections of files/datasets are treated as single entities.
- Aspects of aggregation:
  - combining/merging variables,
  - joining variables,
  - new coordinate axes,
  - overlaying/adding metadata,
  - nesting datasets
- PCMDI CDAT supports aggregations via the cdscan utility that uses XML representation

- cdscan will analyse the archive for:
  - variable information
  - axis information
  - global (universal) metadata

- Why use cdscan
  - Large datasets described as a grouped entity.
  - No need to know underlying data format.
  - No need to know file-names.
  - Datasets can be sliced in any way the user chooses using logical spatio-temporal selectors rather than loops of programming code.
  - You can use it to improve the metadata of your data files…
- cdscan in action
  - $ cdscan –x monthly_means.xml ./*.nc

**MV**

FIGURE 1. Relationships between array abstractions.

NumPy Array — Numeric

NumPy Array — Mask — MA

NumPy Array — Mask — Domain Information — MV
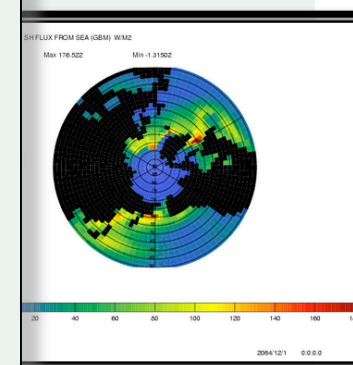+ Metadata

```
>>> import cdms, MV
>>> f_surface = cdms.open('sftlf_ta.nc')
>>> surf = f_surface('sftlf')

# Designate land where "surf" has values
# not equal to 100
>>> land_only = MV.masked_not_equal(surf, 100.)
>>> land_mask = MV.getmask(land_only)

# Now extract a variable from another file
>>> f = cdms.open('ta_1994-1998.nc')
>>> ta = f('ta')

# Apply this mask to retain only land values.
>>> ta_land = cdms.createVariable(ta,
        mask=land_mask, copy=0, id='ta_land')
```
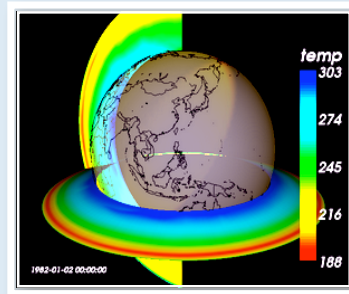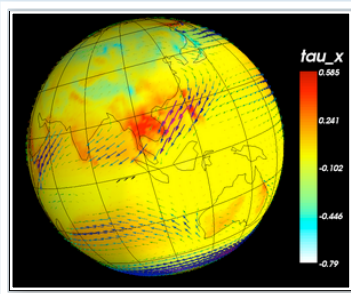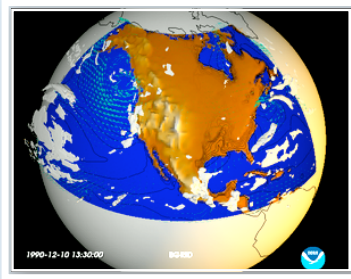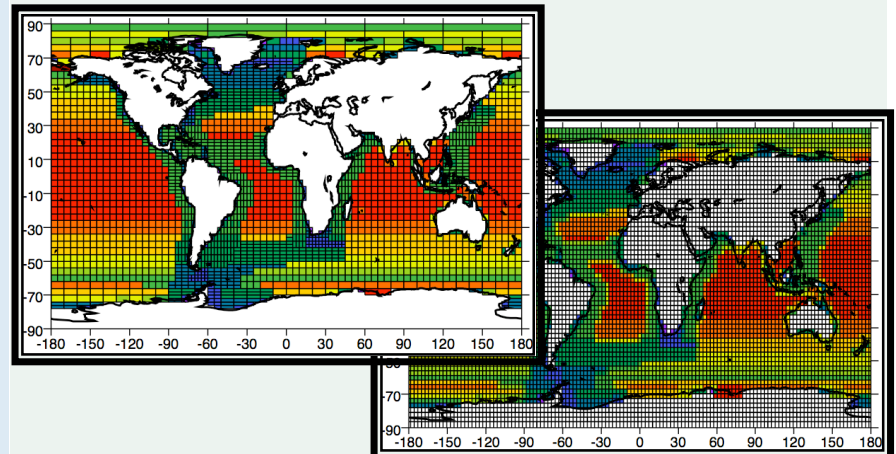
# CDAT examples



## Ncvtk

**Collaboration:**

CDAT developers are currently working with Ncvtk developers to make Ncvtk 3D graphics accessible to the CDAT community. Ncvtk is a collection of commonly used 3D visualization methods applied to data on structured lat/lon grids.



## Regridder

```
#!/usr/local/cdat/bin/python
import cdms
from regrid import Regridder
f = cdms.open('temp.nc')
t= f.variables['t']
ingrid = t.getGrid()
outgrid = cdms.createUniformGrid( -90.0, 46, 4.0, 0.0, 72, 5.0)
regridFunc = Regridder(ingrid, outgrid)
newt = regridFunc(t)
import vcs
vcs.init().plot(t)
vcs.init().plot(newt)
```
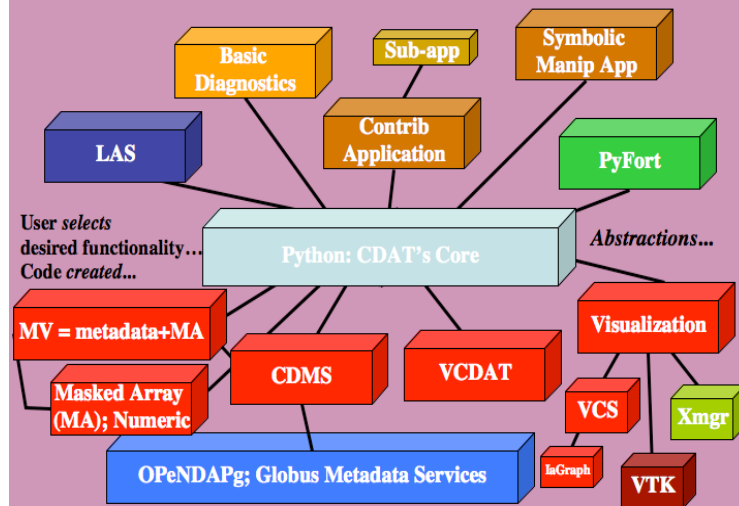
# CDAT facts and figures

## CDAT Users

- **Over 120 mailing list registers**
  - ➤ Probably 7 to 10 times more casual users
- **Mailing list archive: over 1,000 message (~30 per month)**
- **912 Downloads since May 19, 2006 for version 4.1**
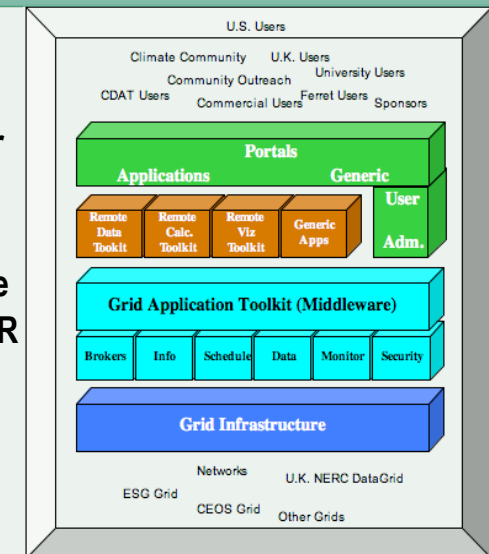- **Improved Documentation**

## CDAT Core Modules



## CDAT Collaborations

**Some CDAT development centers:**

- British Atmospheric Data Center
- LBNL
- GFDL
- Laboratory of Science of Climate and the Environment (LSCE), FR
- PCMDI
- University of Chicago
- University of Hawaii
- University of Reading, UK

# Simple intercomparison use case scenario

| Current Scenario | Future Scenario |
|---|---|
| • Browse PCMDI's centralized database<br>• Download data<br>• Organize data on local site<br>• Regrid data at local site<br>• Perform diagnostics<br>• Produces results | • Search, browse and discover distributed data<br>• Remote site<br>   ➢ Request data<br>   ➢ Regrids<br>   ➢ Diagnostics<br>• ESG returns results |